



PATENT APPLICATION

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

In re the Application of:

Jukka A. VAINIO et al.

Group Art Unit: 2154

Application No.: 10/630,972

Filed: July 31, 2003

Attorney Dkt. No.: 60091.00219

For: DATA COLLECTION IN A COMPUTER CLUSTER

CLAIM FOR PRIORITY UNDER 35 USC § 119

Commissioner for Patents
P.O. Box 1450
Alexandria, VA 22313-1450

November 10, 2003

Sir:

The benefit of the filing dates of the following prior foreign application filed in the following foreign country is hereby requested for the above-identified patent application and the priority provided in 35 U.S.C. §119 is hereby claimed:

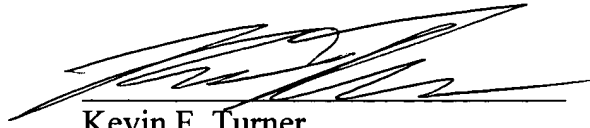
Finnish Patent Application No. 20030796 filed on May 27, 2003 in Finland

In support of this claim, a certified copy of said original foreign application is filed herewith.

It is requested that the file of this application be marked to indicate that the requirements of 35 U.S.C. §119 have been fulfilled and that the Patent and Trademark Office kindly acknowledge receipt of this document.

Please charge any fee deficiency or credit any overpayment with respect to this paper to Counsel's Deposit Account No. 50-2222.

Respectfully submitted,

A handwritten signature in black ink, appearing to read 'Kevin F. Turner', written over a horizontal line.

Kevin F. Turner
Registration No. 43,437

Customer No. 32294
SQUIRE, SANDERS & DEMPSEY LLP
14TH Floor
8000 Towers Crescent Drive
Tysons Corner, Virginia 22182-2700
Telephone: 703-720-7800
Fax: 703-720-7802

KFT:lls

Enclosure: Priority Document (1)

PATENTTI- JA REKISTERIHALLITUS
NATIONAL BOARD OF PATENTS AND REGISTRATION

Helsinki 4.6.2003

ETUOIKEUSTODISTUS
PRIORITY DOCUMENT



Hakija
Applicant

Nokia Corporation
Helsinki

Patenttihakemus nro
Patent application no

20030796

Tekemispäivä
Filing date

27.05.2003

Kansainvälinen luokka
International class

G06F

Keksinnön nimitys
Title of invention

"Data collection in a computer cluster"
(Tiedonkeruu tietokoneklusterissa)

Täten todistetaan, että oheiset asiakirjat ovat tarkkoja jäljennöksiä Patentti- ja rekisterihallitukselle alkuaan annetuista selityksestä, patenttivaatimuksista, tiivistelmästä ja piirustuksista.

This is to certify that the annexed documents are true copies of the description, claims, abstract and drawings originally filed with the Finnish Patent Office.

Eija Solja
Apulaistarkastaja

Maksu 50 €
Fee 50 EUR

Maksu perustuu kauppa- ja teollisuusministeriön antamaan asetukseen 1027/2001 Patentti- ja rekisterihallituksen maksullisista suoritteista muutoksineen.

The fee is based on the Decree with amendments of the Ministry of Trade and Industry No. 1027/2001 concerning the chargeable services of the National Board of Patents and Registration of Finland.

Osoite: Arkadiankatu 6 A
P.O.Box 1160
FIN-00101 Helsinki, FINLAND

Puhelin: 09 6939 500
Telephone: + 358 9 6939 500

Telefax: 09 6939 5328
Telefax: + 358 9 6939 5328

22

1

DATA COLLECTION IN A COMPUTER CLUSTER

Field of the Invention

- 5 [0001] The present invention relates generally to computer clusters that include a plurality of computer nodes. More particularly, the present invention relates to a mechanism for collecting state information within the cluster. In this context, state information refers to data that indicates how the resources of a computer node are able to complete their tasks in the cluster. The state
- 10 information may thus include, not only data indicating the current load of the various resources of a computer node, but also data about the current performance or capacity of the resources in the computer node, i.e. data about the current ability of the resources to complete their tasks in the cluster.

Background of the Invention

- 15 [0002] As is commonly known, a computer cluster is a group of computers working together to complete one or more tasks. Computer clusters can be used for load balancing, for improved fault tolerance (i.e. for improved availability in case of failures), or for parallel computing, for example.
- 20 [0003] A typical computer cluster comprises a plurality of computer nodes. A computer node here refers to an entity provided with a dedicated processor, memory, and operating system, as well as with a network interface through which it can communicate with other computer nodes of the cluster. At least one of the computer nodes in the cluster is capable of acting as a manager
- 25 node that manages the cluster. In order to detect failures in the cluster, the manager node sends certain messages, called heartbeats, periodically to the other computer nodes in the cluster. Typically, only one computer node at a time acts as a manager node.
- 30 [0004] Control software, residing typically in the manager node, has to monitor all computer nodes that belong to the cluster. In order to get a true and up-to-date picture of the state of the nodes, the control software has to collect state information at a fairly high frequency from the nodes. This is a problem especially in large computer clusters, which may contain tens, or even

hundreds of computer nodes. In these large computer clusters the data collection rate has to be compromised in favor of the performance of the network and the computer nodes, to ensure that the network does not become congested due to the data collection and that the performance of the computer nodes remains at an acceptable level despite the data collection performed. In other words, in large clusters the data collection rate has to be compromised in order not to degrade the performance of the network or the computer nodes excessively.

[0005] The objective of the present invention is to eliminate or alleviate this drawback.

Summary of the Invention

[0006] One objective of the invention is to bring about a novel mechanism for collecting state information from the computer nodes of a computer cluster. A further objective of the invention is to bring about a mechanism that does not require the collection rate of the state information to be compromised in favor of network or node performance even in large clusters.

[0007] In the present invention, an internal property of a computer cluster, the heartbeat mechanism, is utilized for collecting state information from the computer nodes for monitoring and control purposes. As described below, the collected state information may be utilized either internally in the computer cluster or by an outside entity, such as a network monitoring or management system.

[0008] Thus one embodiment of the invention is the provision of a method for transferring state information in a computer cluster comprising a plurality of computer nodes. The method includes the steps of:

- transmitting a heartbeat message from a first computer node of a computer cluster to a second computer node of the computer cluster, the second computer node including at least one resource for performing at least one cluster-specific task;
- receiving the heartbeat message in the second computer node;
- retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state

information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and

- sending the state information in the heartbeat acknowledgment message to the first computer node.

5 **[0009]** In a further embodiment the invention provides a computer cluster comprising a plurality of computer nodes. The computer cluster includes:

- first means for transmitting a heartbeat message from a first computer node of the computer cluster to a second computer node of the computer cluster, the second computer node including at least one resource

10 for performing at least one cluster-specific task;

- second means for receiving the heartbeat message in the second computer node;

15 - third means for retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and

- fourth means for sending the state information in the heartbeat acknowledgment message to the first computer node.

20 **[0010]** In another embodiment the invention provides a computer node for a computer cluster. The computer node includes:

- at least one resource for performing at least one cluster-specific task;

- first means for receiving a heartbeat message from another computer node;

25 - second means for retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and

30 - third means, responsive to the second means, for sending the state information in the heartbeat acknowledgment message to said another computer node.

[0011] By means of the solution of the invention, real-time state information can be collected from the computer nodes of a computer cluster without excessively loading the network or the computer nodes, i.e. the information

collection rate does not need to be compromised due to the load the collection causes. The overhead caused by the increased length of the acknowledgment message is relatively low, especially if the length of the minimum transmission unit is not exceeded.

- 5 **[0012]** In one embodiment of the invention, a computer node receiving a heartbeat message checks whether state information is to be retrieved for the heartbeat acknowledgment message to be sent as a response to the heartbeat message. In this way, unnecessary transfer of state information can be avoided.
- 10 **[0013]** A further advantage of the invention is that the collected information may be simultaneously utilized by different entities within or outside of the computer cluster.
- 15 **[0014]** Other features and advantages of the invention will become apparent through reference to the following detailed description and accompanying drawings.

Brief Description of the Drawings

- 20 **[0015]** In the following, the invention and its preferred embodiments are described more closely with reference to the examples shown in FIG. 1 to 5 in the appended drawings, wherein:
- [0016]** FIG. 1 illustrates one computer cluster according to the invention;
- [0017]** FIG. 2 is a flow diagram illustrating the basic operation of a manager node in view of one heartbeat message;
- 25 **[0018]** FIG. 3a is a flow diagram illustrating one embodiment for sending state information from a computer node;
- [0019]** FIG. 3b is a flow diagram illustrating another embodiment for sending state information from a computer node;
- [0020]** FIG. 4 is a schematic diagram illustrating the collection of state information in a computer node; and
- 30 **[0021]** FIG. 5 is a schematic presentation of a heartbeat message according to the invention.

Detailed Description of the Invention

[0022] FIG. 1 shows an example of a computer cluster **100** in which the mechanism of the invention is utilized. The cluster comprises N computer nodes **110_i** ($i=1,2,3,\dots,N$). Each computer node is an independent entity provided with a processor, memory and an operating system copy of its own. Each computer node is further provided with a network interface for connecting it to a network **120**, which is typically an Internet Protocol (IP) based network. It is to be noted here that the mechanism of the invention is not dependent on the transmission protocol, but may be applied in many different environments. However, an IP network forms a typical environment for the invention.

[0023] At each time, one of the computer nodes, in this example node **110₁**, operates as a manager node that manages the cluster and its resources. In order to detect failures occurring in the cluster, the manager node sends heartbeat messages **HB** periodically to the other computer nodes in the cluster. Although the cluster may include more than one node being able to act as a manager node, one of such nodes operates as the manager node at a time. A single heartbeat message is typically a multicast message destined for all nodes of the cluster, and the period between two successive heartbeat messages depends greatly on the application environment.

[0024] When a computer node receives a heartbeat message from the manager node, it returns a heartbeat acknowledgment message **HB_ACK** to the manager node, indicating to the manager node that it is alive and can therefore remain in the cluster. If the manager node does not receive a heartbeat acknowledgment message from a computer node, it starts recovery measures immediately. Typically, the computer node with which a communication failure has been detected is removed from the cluster, and the cluster-specific activities of the node are reassigned to one or more other nodes.

[0025] A variety of different tasks may be performed by the cluster, and the actual applications may be distributed in a variety of ways within the cluster. One or more of the cluster nodes may appear as a single entity to an element external to the cluster. For example, if the computer nodes perform routing, one or more of the computer nodes may form a routing network element, as

seen from the outside of the cluster. In an extreme case, all computer nodes appear as a single entity to an external viewer.

5 [0026] If load sharing groups are utilized in the cluster, one or more of the computer nodes may further operate as an Internet Protocol Director (IPD) node, which is a load sharing control node routing incoming task requests within a load sharing group. In the example of FIG. 1, computer node 110₂ operates as an IPD node receiving task requests from the outside of the computer cluster.

10 [0027] In the present invention, the intrinsic heartbeat mechanism of a computer cluster is utilized for collecting state information from the computer nodes. The data may be collected for the purposes of the cluster only, or for an entity external to the cluster, such as a network monitoring or management system 160 connected to the network. The heartbeat acknowledgment messages are used to carry state information from the cluster nodes to the
15 manager node, which then stores the information in a Management Information Base (MIB) 150.

20 [0028] In one embodiment of the invention, the MIB is made available for both entities within the computer cluster and for entities external to the computer cluster. For example, the internal fault management of the cluster may utilize the data collected. The fault management logic may be distributed in the cluster with an agent 130 residing in the manager node so that the fault management system can read data from the MIB. In other words, the fault management system may comprise a client-server mechanism with the server part residing in the manager node and the client parts residing in the computer
25 nodes. Another cluster entity capable of utilizing the MIB is a computer node that allocates incoming tasks to the computer nodes performing said tasks. In addition to the above-mentioned IPD node, any other cluster node may operate as such a load balancing entity.

30 [0029] Access to the MIB can be implemented in any known manner either directly or through the manager node, depending on whether the MIB forms an independent network node or whether it is connected to the manager node. The MIB may also be connected to a computer node other than the manager node.

[0030] FIG. 2 is a flow diagram illustrating an example of the basic operation of the manager node with respect to one heartbeat message sent to another computer node. It is thus to be noted here that FIG. 2 illustrates the operation with respect to one heartbeat message sent, i.e. the periodic sending of the heartbeat messages is not shown in the figure. When the manager node transmits a heartbeat message, it sets a timer (step 201) and starts to monitor if a heartbeat acknowledgment message is received as a response from said another computer node (step 202). If this acknowledgment message arrives before the expiration of the timer, the manager node examines the message (step 204). If the manager node detects the message contains state information, it extracts the said information from the message and updates the MIB based on the information (step 207). In case of an acknowledgment message void of state information the manager node proceeds in a conventional manner.

[0031] If the timer expires before a heartbeat acknowledgment message is received, the manager node concludes that a communication failure has occurred with the computer node, and starts recovery measures (step 205). In practice, the time period measured by the timer is so long that more than one heartbeat messages can be transmitted within that period. A heartbeat acknowledgment received for any of these messages then triggers the process to jump to step 204. Normally the manager node proclaims a computer node to be faulty when N successive heartbeat messages remain without an acknowledgment from that computer node. The manager node may thus be allowed to lose a given number of heartbeat messages before the recovery measures are started. Particularly in case of the UDP (User Datagram Protocol), which is commonly used for carrying heartbeat messages, messages may be lost without a real problem existing in the network. In view of the above, FIG. 2 is to be seen merely as an illustration of the processing principles of the incoming heartbeat acknowledgment messages in the manager node, while the actual implementation of the relevant manager node algorithm may vary in many ways.

[0032] FIG. 3a is a flow diagram illustrating an example of the operation of a computer node with respect to one heartbeat message received from the manager node. When the heartbeat message is received, the computer node

- examines (step 301) whether a predetermined condition is fulfilled. This predetermined condition is set in order not to transfer state information unnecessarily in the acknowledgment messages. If the condition is fulfilled, the computer node retrieves state information from its memory (step 303) and
- 5 generates a heartbeat acknowledgment message containing the state information retrieved. If the predetermined condition is not fulfilled, the computer node generates a normal heartbeat acknowledgment message, i.e. a heartbeat acknowledgment message without state information (302). The generated message is then sent back to the manager node (step 305).
- 10 **[0033]** The predetermined condition set for the retrieval of the state information is typically such that a certain minimum time period must have passed since the latest transmission of state information to the manager node. If this time limit has been exceeded, new state information is retrieved and inserted into the heartbeat acknowledgment message. Otherwise a normal
- 15 heartbeat acknowledgment message is sent. In order to detect when the time limit has been exceeded, the computer node may start a counter at step 305. The current value of the counter is then examined at step 301 in connection with a subsequent heartbeat message. The computer node thus typically sends both normal heartbeat acknowledgment messages and heartbeat
- 20 acknowledgment messages containing the state information, the proportions of these two message types depending on the rate of the heartbeat messages received.
- [0034]** The predetermined condition set for the retrieval of the state information may also consist of several sub-conditions that must be fulfilled
- 25 before state information is retrieved. If the load of the computer node is used as such a sub-condition, the retrieval of the state information could occur, for example, only if both a certain minimum time period has passed since the latest transmission of state information and the current load of the computer node is below a certain maximum level.
- 30 **[0035]** As shown in FIG. 3b, it is also possible that the node determines, in response to the reception of a heartbeat message, the type of state information to be retrieved (step 311). Different types of information may thus be carried by successive heartbeat acknowledgment messages. For example, if heartbeat messages are transmitted frequently enough, a certain set of

parameters may be carried by N successive heartbeat acknowledgment messages, the same set being again transmitted by the next N heartbeat acknowledgment messages, and so on. Furthermore, certain information (parameters) may be transferred less frequently than other information.

- 5 **[0036]** The state information retrieved from the memory depends generally on the application running on the computer node. However, certain basic parameters that relate to the operating system of the computer node are the same for all computer nodes. These parameters include figures indicating the CPU idle time and the number of certain I/O operations, for example.
- 10 Basically, the state information can be divided into two groups: the parameters relating to the performance of the applications and the parameters relating to the performance and/or state of the node platform.

- [0037]** FIG. 4 illustrates an example of the software architecture of the heartbeat acknowledgment generation in a computer node. A kernel module
- 15 **400** residing in the kernel space receives the parameters relating to the operating system directly from the kernel space of the computer node. In the user-space, where the applications are executed, each application **401** may have a library **402** through which it can write the relevant parameters to the kernel module. A supervision agent **403** residing in the user space retrieves
- 20 the state information from the kernel module if the predetermined condition is fulfilled, and constructs the heartbeat acknowledgment message containing the information retrieved. In the embodiment of FIG. 4, the storage of the state information is thus implemented in the operating system, which provides a faster operation. However, the state information may also be stored in a mass
- 25 memory, such as a disk.

- [0038]** FIG. 5 illustrates a general structure of the heartbeat acknowledgment message containing state information. The message comprises three successive portions: a header portion **501** that includes the protocol headers of the relevant protocols (such as Ethernet, IP and TCP/UDP headers), an acknowledgment identifier **502**, and a payload portion **503** that contains the
- 30 state information retrieved in the computer node. The message is thus otherwise similar to a conventional heartbeat acknowledgment message, but it includes a payload portion that contains the state information. In one embodiment of the invention the payload portion is encoded by using ASN.1

(Abstract Syntax Notation One) and PER (Packed Encoding Rules) coding. In this way the state information can be packed efficiently and more information can be inserted into the same message space. Depending on the protocols used, part of the state information may be transmitted without causing any
 5 extra load in the network. This is the case if the length of a conventional heartbeat message is shorter than the length of the minimum transmission unit, in which case state information may be used as the padding bits.

[0039] The load increase caused by a heartbeat acknowledgment message of the invention is relatively small as compared to the load caused by a
 10 conventional heartbeat acknowledgment message. This is because the overhead caused by a longer message is relatively low, since in short messages the protocol header takes a major part of the transmitted message. Furthermore, as messages shorter than a minimum message length are normally filled up, they may now be filled with the state information. In this way
 15 part of the state information may be transferred without causing extra load in the network. The extra load caused by the method of the invention therefore also depends on the environment where the invention is applied. In an Ethernet network, for example, this minimum message length is 64 bytes, which is more than portions **501** and **502** require.

[0040] Although the invention was described above with reference to the
 20 examples shown in the appended drawings, it is obvious that the invention is not limited to these, but may be modified by those skilled in the art without departing from the scope and spirit of the invention. For example, it is not necessary to check whether a normal heartbeat acknowledgment message or
 25 a heartbeat acknowledgment message containing state information is to be sent, but an acknowledgment message containing state information can be sent in response to every heartbeat message.

Claims

83

1. A method for transferring state information in a computer cluster comprising a plurality of computer nodes, the method comprising the steps of:
 - 5 - transmitting a heartbeat message from a first computer node of a computer cluster to a second computer node of the computer cluster, the second computer node including at least one resource for performing at least one cluster-specific task;
 - receiving the heartbeat message in the second computer node;
 - 10 - retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and
 - sending the state information in the heartbeat acknowledgment message to the first computer node.
2. A method according to claim 1, further comprising the step of examining, in response to the receiving step, whether state information is to be retrieved for the heartbeat acknowledgment message.
3. A method according to claim 2, wherein the examining step
 - 15 includes examining whether a predetermined condition is fulfilled.
4. A method according to claim 3, wherein the retrieving and sending steps are performed when the examining step indicates that the predetermined condition is fulfilled, and wherein the method further comprises the step of
 - 20 transmitting a heartbeat acknowledgment message without state information when the examining step indicates that the predetermined condition fails to be fulfilled.
5. A method according to claim 1, further comprising the step of determining the type of state information to be retrieved for the heartbeat acknowledgment message.
6. A method according to claim 1, further comprising the step of
 - 25 storing the state information sent to the first computer node in a Management Information Base (MIB).
7. A method according to claim 6, further comprising the step of
 - 30 transferring data from the Management Information Base to an entity external to the computer cluster.
 - 35

8. A computer cluster comprising a plurality of computer nodes, the computer cluster comprising:

5 - first means for transmitting a heartbeat message from a first computer node of the computer cluster to a second computer node of the computer cluster, the second computer node including at least one resource for performing at least one cluster-specific task;

- second means for receiving the heartbeat message in the second computer node;

10 - third means for retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and

- fourth means for sending the state information in the heartbeat acknowledgment message to the first computer node.

15 9. A computer cluster according to claim 8, further comprising a Management Information Base (MIB) operably connected to the first computer node for storing the state information sent to the first computer node.

20 10. A computer cluster according to claim 9, further comprising first access means for accessing the Management Information Base from the computer cluster.

11. A computer cluster according to claim 9, further comprising second access means for accessing the Management Information Base from outside of the computer cluster.

25 12. A computer cluster according to claim 11, wherein the second access means comprise a network interface in the first computer node.

13. A computer node for a computer cluster, the computer node comprising:

- at least one resource for performing at least one cluster-specific task;

30 - first means for receiving a heartbeat message from another computer node;

35 - second means for retrieving state information for a heartbeat acknowledgment message to be sent as a response to said heartbeat message, the state information indicating the ability of said at least one resource to perform said at least one cluster-specific task; and

- third means, responsive to the second means, for sending the state information in the heartbeat acknowledgment message to said another computer node.

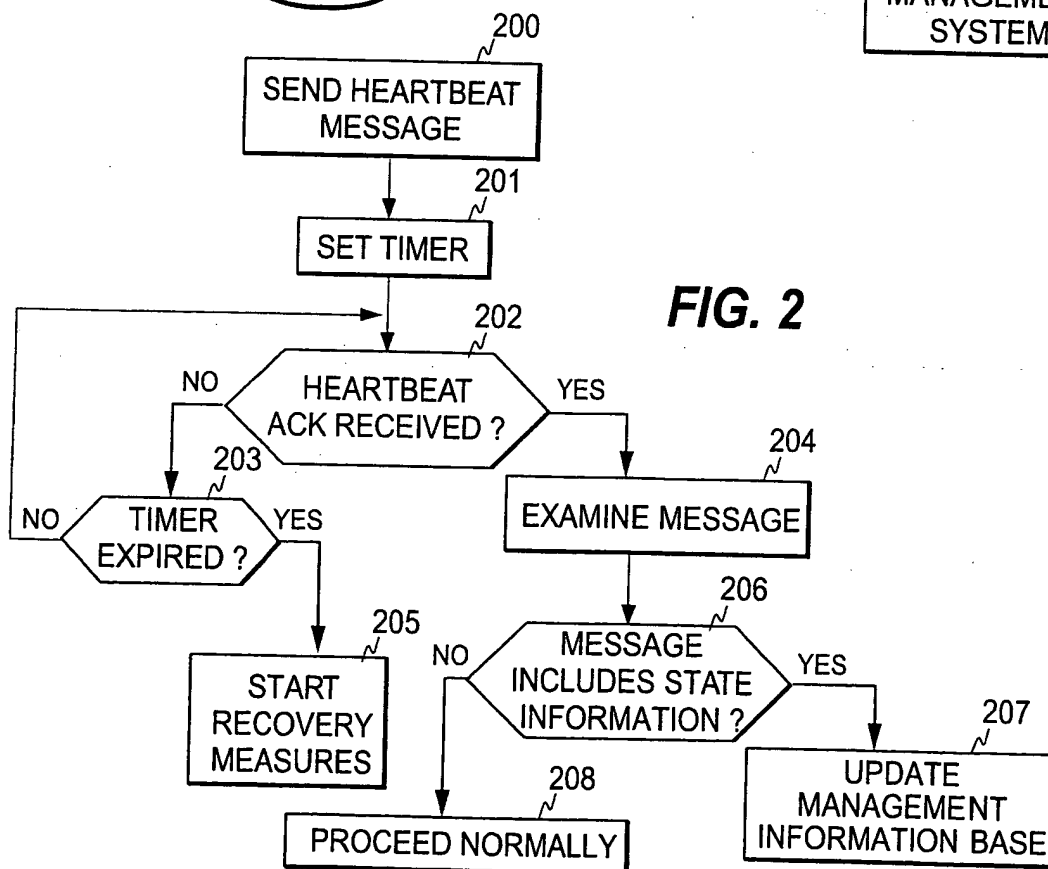
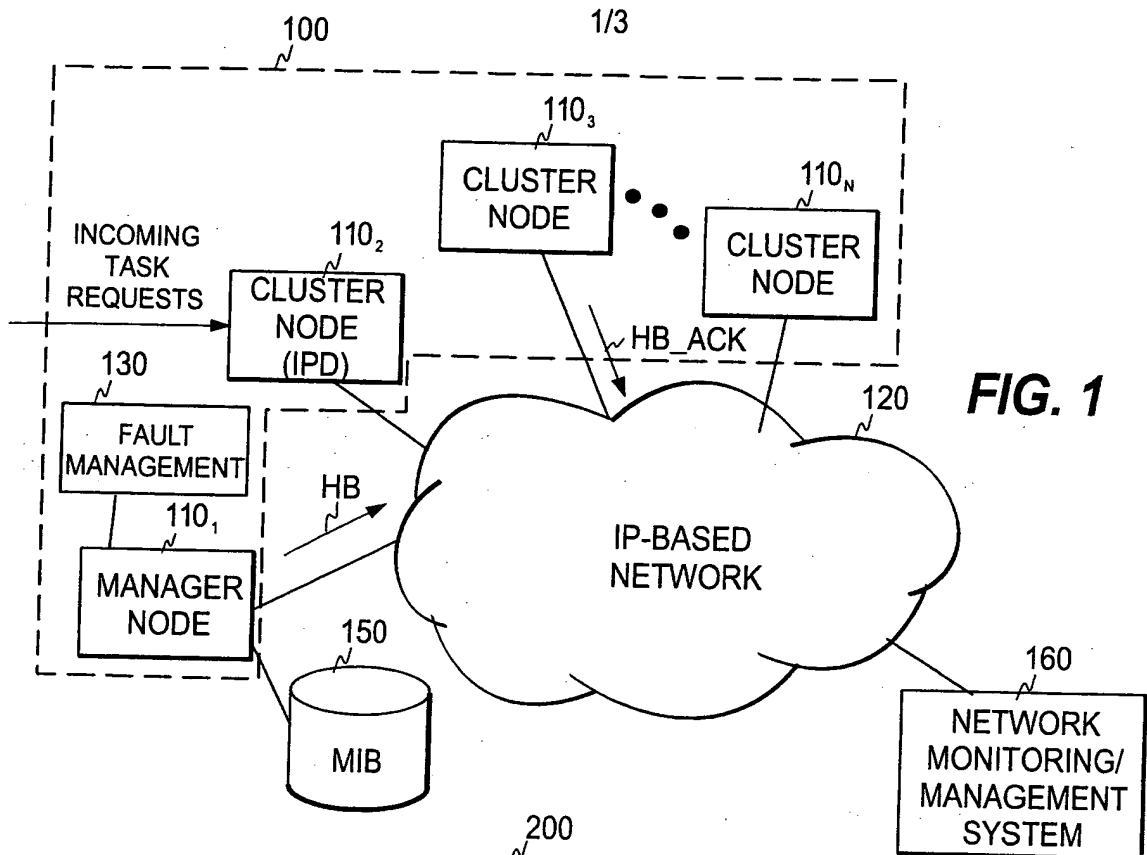
- 5 14. A computer node according to claim 13, further comprising fourth means for examining whether state information is to be retrieved for the heartbeat acknowledgment message.

L4

(57) Abstract

The invention relates to a mechanism for transferring state information in a computer cluster comprising a plurality of computer nodes. In the method, heartbeat messages are sent periodically from a first computer node of the computer cluster to other computer nodes of the cluster. Each of said other nodes includes at least one resource for performing at least one cluster-specific task. In order that up-to-date state information could be collected even in large clusters about the ability of the resources to perform the cluster-specific tasks, without excessively loading the computer nodes and the network, current state information is returned in a heartbeat acknowledgment message to the node that sent the heartbeat message.

45



25

2/3

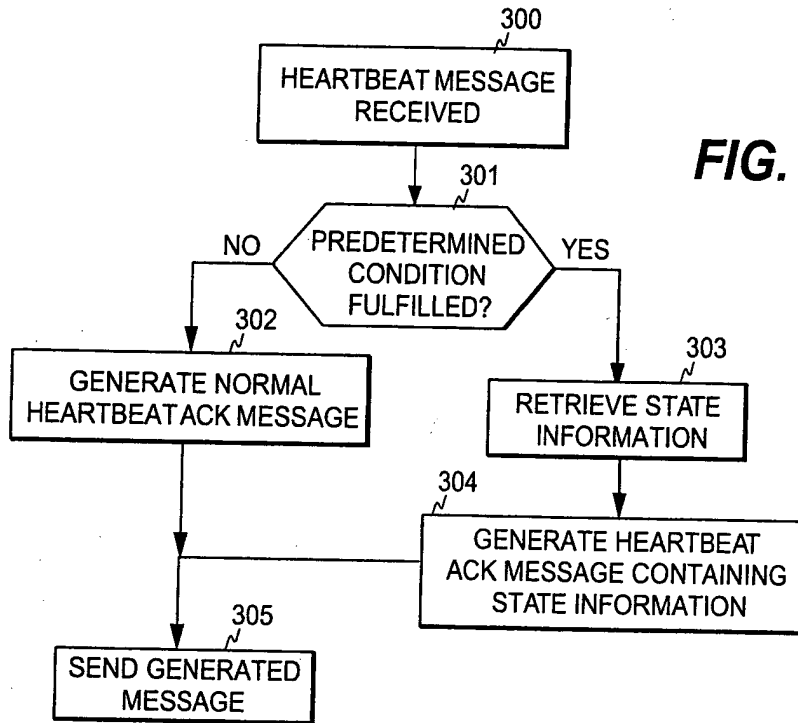


FIG. 3a

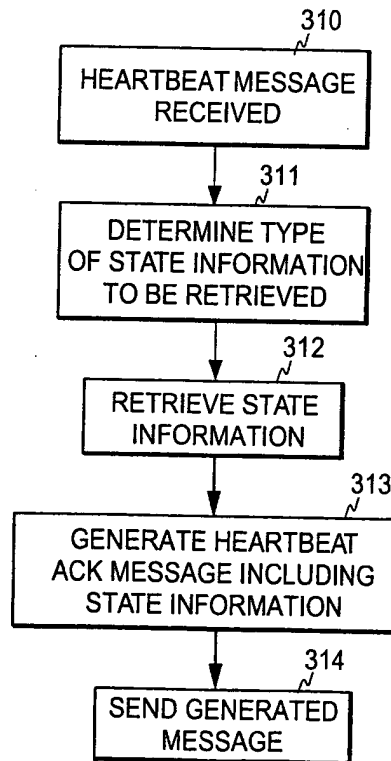


FIG. 3b

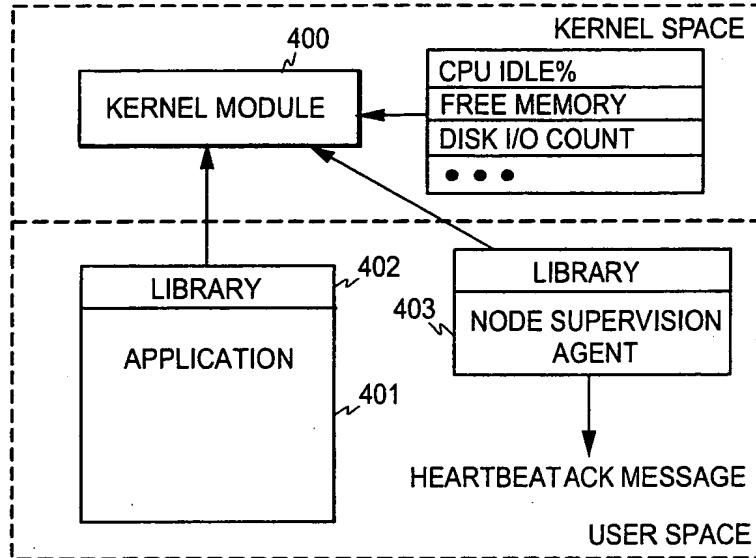


FIG. 4

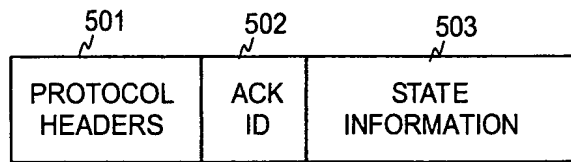


FIG. 5